

## CLAIMS

What is claimed is:

1 1. A method of categorizing a plurality of new electronic documents into a set of  
2 categories, each of the categories containing a plurality of training set documents, by  
3 using a matrix representing document similarity that is derived by combining two or  
4 more measures of document similarity.

1 2. A method as recited in Claim 1, wherein the measures of document similarity include  
2 hyperlink similarity.

1 3. A method as recited in Claim 2, in which two documents among the plurality of  
2 documents are considered similar to each other when there is a link from one to the  
3 other, or when the two documents link to, or are linked to by, a set of other associated  
4 documents.

1 4. A method as recited in Claim 3, in which certain hyperlinks have greater or lesser  
2 similarity weight than other hyperlinks, based on other features of the links or their  
3 source or destination documents.

1 5. A method as recited in Claim 1, wherein the measures of document similarity include  
2 a similarity of text of the documents.

1 6. A method as recited in Claim 5, wherein two documents are considered similar based  
2 on a comparison of word vectors derived from the text of each of the two documents.

1 7. A method as recited in Claim 5, wherein text similarity is determined in part based  
2 upon weight values assigned to words of the text, and wherein certain words have  
3 greater or lesser weight than other words.

1 8. A method as recited in Claim 1, wherein the measures of document similarity include  
2 user click-through similarity.

1 9. A method as recited in Claim 8, wherein two documents are considered similar based  
2 on user click-through similarity when the documents are associated with similar  
3 patterns of user click behavior, selected from among frequency of clicks, click  
4 context, duration of viewing, proximity in time to other clicks, or proximity in context  
5 to other clicks.

1 10. A method as recited in Claim 1, wherein the measures of document similarity are  
2 derived from patterns detected in user viewing of the documents.

1 11. A method as recited in Claim 10, wherein the user viewing information is monitored  
2 by a web caching system and stored in a log.

1 12. A method as recited in Claim 10, wherein two documents are considered similar based  
2 on patterns of user viewing behavior, including frequency of viewing, viewing  
3 context, duration of viewing, proximity in time to other documents viewed by the  
4 same user, or similarity of patterns of viewing by all users.

1 13. A method as recited in Claim 1, wherein the measures of document similarity include  
2 URL similarity.

1 14. A method as recited in Claim 13, wherein two documents are considered similar if a  
2 URL of each document contains similar URL sub-components.

1 15. A method as recited in Claim 1, wherein the measures of document similarity include  
2 multimedia similarity.

1 16. A method as recited in Claim 15, wherein two documents are considered similar based  
2 on features derived from multimedia components linked to or contained by the  
3 documents.

1 17. A method as recited in Claim 1, wherein the combination of two or more measures of  
2 document similarity is achieved by taking the union of each of a plurality of graphs,  
3 each graph describing one of the measures of document similarity, to compute a  
4 combined graph that describes the combined document similarity.

1 18. A method as recited in Claim 1, wherein the combination of two or more measures of  
2 document similarity is achieved by taking the intersection of each of a plurality of  
3 graphs, each graph describing one of the measures of document similarity, to compute  
4 a combined graph that describes the combined document similarity.

4/5  
7/9

2

3

1

2

1

2

1

2

1

2

1

2

3

4

1

2

19. A method as recited in Claim 1, further comprising the step of extracting structural information from the similarity matrix to obtain new documents supported by the set of training documents for each category.

20. A method as recited in Claim 19, wherein the structural information is obtained by optimizing an objective function.

21. A method as recited in Claim 19, wherein the structural information is obtained by only approximately optimizing an objective function.

22. A method as recited in Claim 21, wherein approximately optimizing the objective function comprises repeated application of a growth transformation.

23. A method as recited in Claim 19, further comprising the step of creating and storing a second matrix that represents an interim score for each document in each category.

24. A method as recited in Claim 19, further comprising the steps of, periodically as the matrix is being computed, normalizing rows of the matrix by normalizing within each document, across all categories, whereby the score for one document in a particular category will depend on the scores for that document in all other categories.

25. A method as recited in Claim 19, further comprising the steps of, periodically as the matrix is being computed, normalizing columns of the matrix by normalizing within

3 each category, across all documents, whereby the score for one document in a  
4 particular category depends on the scores for all other documents in that category.

1 26. A method as recited in Claim 1, in which the categories come from a manually  
2 defined taxonomy.

1 27. A method as recited in Claim 1, wherein the categories are derived from logs of user  
2 queries.

28. A method as recited in Claim 1, further comprising the steps of creating and storing  
the matrix using columns representing documents and rows representing user  
sessions, and wherein values of elements of the second matrix represent interest in a  
document shown by a particular user in a particular session.

29. A method as recited in Claim 1, further comprising the steps of creating and storing  
the matrix using columns representing user sessions and rows representing  
documents, and wherein values of elements of the second matrix represent interest in  
a document shown by a particular user in a particular session.

30. A method as recited in Claim 28, wherein the element values are computed as a  
function of a time that a user has spent viewing a document associated with each  
element.

1 <sup>30</sup>  
31. A method as recited in Claim 28, further comprising the steps of creating and storing a  
2 second matrix representing a Similarity between pairs of documents i and j, wherein  
3 the second matrix is derived by comparing pairs of column vectors or row vectors,  
4 respectively i and j of the first matrix.

1 <sup>31</sup>  
32. A method as recited in Claim 28, further comprising the steps of creating and storing a  
2 second matrix representing a Similarity between pairs of documents i and j, by finding  
3 pairs of documents i and j which have high interest values for a particular user in a  
4 particular session or period of time.

1 33. The method recited in Claim 1, further comprising the steps of:  
2 identifying a category of a classification taxonomy of the hypertext system in which a  
3 first electronic document is presently classified; and  
4 if a second electronic document is found to be highly Similar, storing information that  
5 classifies the second electronic document into the category.

1 <sup>34</sup>  
34. A computer-readable medium carrying one or more sequences of instructions, wherein  
2 execution of the one or more sequences of instructions by one or more processors  
3 causes the one or more processors to perform the steps of categorizing a plurality of  
4 new electronic documents into a set of categories, each of the categories containing a  
5 plurality of training set documents, by using a matrix representing document  
6 similarity that is derived by combining two or more measures of document similarity.

1 35. A method of categorizing a plurality of new electronic documents for use in a  
2 hypertext search system, the method comprising the steps of:  
3 creating and storing a set of categories for the documents;  
4 creating and storing a matrix, in which rows and columns identify documents, and in  
5 which each element of the matrix stores a value that represents a similarity  
6 among a pair of documents associated with a row and column that intersect at  
7 the element;  
8 deriving each matrix value by combining two or more measures of similarity that are  
9 obtained by analysis of the documents.

1 36. A method as recited in Claim 35, further comprising the steps of:  
2 for each measure of document similarity, creating and storing a graph of links;  
3 creating and storing a combined graph that combines the graphs and that represents a  
4 generalized similarity of the documents;  
5 computing a generalized similarity value for a pair of documents based on the  
6 combined graph.

1 37. A method as recited in Claim 36, further comprising the steps of classifying  
2 unclassified documents into category nodes of a taxonomy structure associated with  
3 the hypertext search system based on the generalized similarity value in combination  
4 with a comparison of a set of pre-classified training set of documents with a set of  
5 unclassified documents, to carry out classification.